## **Global variation in the HIV-1 V3 region**

### John C. Blouin, Esther A. Guzman and Brian T. Foley

MS K710, Los Alamos National Laboratory, Los Alamos, NM 87545

#### Introduction

Due to the immunogenicity and functional importance of the V3 loop, there has been a great deal of interest in the V3 region of the envelope protein, resulting in a large international effort to obtain V3 region sequences. This section, which includes sequences taken from 1651 individuals and complete references, provides an overview of the variation of sequences that span this region.

### Sequences

To best summarize the spectrum of international HIV-1 variants, only one representative viral sequence was included per infected individual. A complete set of references accompanies the sequence alignments, and nomenclature was preserved from the original papers so individuals and isolates can be clearly identified. HIV1 was deleted from the sequence names in this section, as all sequences included here are HIV1. Included with the references when available are brief descriptions of critical features of the sequences. This includes the health status of the individual from whom the virus was derived, whether or not the virus was cultured, and the year the blood sample was taken.

All sequences are prefaced by a subtype association (see phylogenetic clustering below) and a two letter country code to identify the country that the individual resided in at the time that the blood sample was taken. If the person was a recent immigrant and this information was available, we included the country of origin in the references. The two letter code was developed for Internet (Copyright 1992, Lawrence H. Landwater and the Internet Society), and incorporated here based on a suggestion made by Dr. Francine McCutchan. The key to the country codes follows this introduction. Note that this key has been updated for 1996, with several country codes for eastern european nations added.

Sometimes only one viral sequence was available from a person: a clone from an isolate, or a direct sequence of PCR amplified peripheral blood DNA. For other individuals, up to 80 viral sequences from PCR amplified DNA or RNA from blood samples were available. Consequently, over 8000 sequences are represented by the 1651 included in this section. When two sequences were available from a person, one of the two was randomly selected. When three or more sequences were generated from a person, all available sequences were aligned (without regard to different time points of sampling) and either one representative sequence was chosen, or a consensus of the most common base found in each position in the alignment was generated. If there was a tie (e.g., 10 A's, 10 T's), the top base or amino acid in the alignment was used. If a set of sequences from two or more individuals was epidemiologically linked, and genetically very similar, only one sequence from the set was included, preferably the most recently infected. In the sequence description and references section, the short hand "PCR-direct, peripheral blood DNA" is used to signify that viral DNA was amplified from PBMCs, without culturing, and a single "direct" sequence was obtained from the amplification reaction products. The short hand "Consensus, PCR-clones, peripheral blood DNA" signifies that viral DNA was amplified from PBMCs and a set of clones was generated and sequenced from the PCR amplification products. The cloned sequences were aligned and a consensus was generated. In a handful of cases, a particular gp160 clone from an isolate was shown to be expressed and functional using a vaccinia virus T7 expression system. In these cases, the clone rather than the consensus of all sequences from a particular individual is included.

### **Phylogenetic clustering**

Sequences have been organized according to the phylogenetic subtype association (A-J) of their envelope V3 regions only. The original sequence subtype (A-H) designations were defined based on the phylogenetic relationships determined by using both gag and env genes (when possible), are

#### Introduction

approximately genetically equidistant in envelope, and have multiple members. The phylogenetic subtype designations and associations have generally been adopted by the HIV research community, and are now often presented with the publication of new sequences. We have either determined the subtype designations here, if not specified in the original manuscript, or else confirmed the subtype designations of the original manuscripts, and then used the subtypes to organize this section. Generally, confirmations were done by aligning a set HIV-1 V3 region sequences with longer env gene sequences (Part IIIC) that have clear subtype associations, and then using parsimony or neighbor joining trees to determine associations. Some of the shorter gene fragments from this region were given a subtype designation based on Hamming distances, using the similarity function of the MASE program (Faulkner DV, and Jurka J. TIBS 13:321-322 (1988)); these sequences have ".sh" appended to their name to indicate that they were too short for phylogenetic analysis. Parsimony trees were generated using PAUP (David Swofford, Illinois Natural History Survey), and neighbor-joining trees were generated with Kimura distances and a transition to transversion ratio of 1.3 using PHYLIP (Joseph Felsenstein, University of Washington). All available nucleotide sequence information was used for phylogenetic analysis; longer protein sequences were trimmed to be approximately the same length as the majority of the PCR fragments in this region, for the purposes of presentation. Some sequences were difficult to classify, and are included in the "U", or unclassified, section. In addition, it has recently been noted that recombination between HIV-1 occurs when an individual is infected with more than one strain. A meeting was held in Santa Fe, New Mexico in October, 1995 to discuss the implications of recombination and methods for detecting recombinant sequences. Because inter-subtype as well as intra-subtype recombination is known to occur, the subtype designations reported in this section should be interpretted only as pertaining to the V3 region of the envelope gene. For example HIV-1 MAL from Zaire, is known to be recombinant between subtypes A and D, with the V3 loop of env resembling subtype D. D ZR-MAL is still listed with other subtype D sequences in this study, but may be moved to the U (uncertain) group in the future.

The set of sequences used to help resolve subtype associations included at least two sequences from each subtype (A-H), plus a simian immunodeficiency virus outgroup sequence. The sequences were selected based on being "typical" of the subtype they represent based on phylogenetic analysis. The set has changed as more sequences have accumulated. Thus not all subtype designations were based on the same reference set.

### Limitations of phylogenetic analyses

Most of the PCR derived sequences contain a sub-optimal length for phylogenetic analyses, given the level of variability in this region – typically on the order of 250 to 300 nucleotides. Due to this limitation, some of the classifications in this section are uncertain and are our best estimate given the available information. Control studies were performed to compare the phylogenetic clustering of the V3 region using available longer sequences, however, and these studies indicate that our subtype designations based on the V3 region are generally reliable. For 146 sequences, we had an approximately 700 base region of env available representing all of the subtypes A-H. (The limitation in length was due to including the H subtype sequences, which did not cover all of gp120.) After removing positions in the alignment which included gaps, 519 bases were left. When a 298 base V3 region fragment was excised from this set, and neighbor joining trees were constructed using both the 519 base and 298 base long sequences, the phylogenetic subtype designations were consistent in each case. Further, when a subset of longer gp120 sequences was analyzed (92 of the 146), including 935 bases after removing positions in the alignment which included gaps, the subtype designations were again clear in neighbor-joining trees. This indicates that the limited V3 region PCR fragments, which include more than the V3-loop, are generally able to serve as a reliable basis for subtype determination.

Without detailed analysis, genetic recombination between subtypes may obscure phylogenetic relationships between sequences. A characteristic of recombination is an indeterminate place in phylogenetic analyses, and some of the "Uncertain" category sequences may prove to be recombinant genomes upon further inspection. Also, while a subtype designation based on a gene or gene fragment may be correct, recombination events outside the region examined may have occurred. Therefore, care should be taken to not overinterpret the subtype designations. If one is to discuss the subtype

designations of viral isolates based on the data presented here, they should be refer to the designation as "B-like over V3 loop region," rather than as "subtype B".

### Limitations of V3 amino acid consensus sequences

The V3 amino acid consensus sequences generated for each subtype have interesting features; however, one should be wary about assuming that any of the consensus sequences may broadly represent their subtype. Certainly many V3 loop variants in each of the subtypes are extremely divergent from the consensus sequences. These divergent forms may have very different biological and immunological characteristics from viruses which are similar to the consensus. Additionally, because of the relatively small sample size of most of the subtypes, consensus sequences can be dominated by a small group of highly similar sequences, which may in turn be a sampling artifact. Hence, these consensus sequences are "evolving" as new sequences from each subtype become available.

	V3 LOOP					
***	^ *^ ^^^	* ^^^	^^ ^	^^ ^	*	*
ALL_CONSENSUS IIIRSENITDNAKTIIVQLNESVEI	N CTRPNNNTRK.SIHIGPGQAFYATGDI.I.GDIRÇ	AHC NISRTEWNN.TLQQVAKKLREHF	NKT.IIF	PS.SGGDL.EITTHSF	.NCGG.EF	FYCNT
A_CONSENSUS_96 VMNV-P-K-	VR	·VKTQKY		.AN		
B_CONSENSUS_96 VVFF	TERTE	LAKK-IVQ-G.	V-	.NQPVM		
C_CONSENSUS_96L-N-VH	VRTT	KDKER-GAP.	K-	. A	R	
D_CONSENSUS_96L-NIT-	YQ.RTL-TL-T	GDLL	T	.KP		
E_CONSENSUS_96L-NHK	STTVRDK	LY- E-NG-KE.V-KTEKN.		.Q-PM-H-	R	
F_CONSENSUS_96Q-S-THT-Q-	DK	V-G-QR-RAKSP.	AK-	.NSM	R	
G_CONSENSUS_96LVK-I	TF	VKE.MN-TAQKKI	NT-	.NSA	R	
H_CONSENSUS_96 VKT-NKSP-P-	K	LYT-ED-KRHE-VQQN.	Q	.EMM-T-	R	
I_CONSENSUS_96 VVKT-NAKA-K-	E	GNDDKVISEEKRL-P.	K-	.A-P		
J_CONSENSUS_96 VKKT-T-	V -VAGVLE	GRSR-VAY	TQ-	.K-ANPMT-		
O_CONSENSUS_96 -R-MGKS-SG-NTSTINM	T -EG-Q-VQ.E-KMAWYSMGLAAGNGNNS-A	AY- TYNA-D-EK.AKQTAEYLELV	NSNVT	.TMF.NRSSGGDAEVTHL	HFNCHG	
U_CONSENSUS_96NQ-	TRRT	·		.KP		

**Subtype consensus sequences.** This V3 region alignment shows a consensus sequence generated for each of the eleven subtypes. The subtype consensus sequences indicate the most common amino acid found in each position among the sequences associated with each subtype. The sequences are aligned to a consensus based on the most common amino acid in the subtype consensus sequences, which approximates a "global" consensus. It was generated in this way (rather than by using all 1651 sequences) to avoid over-representation of the B subtype, which has by far the largest number of available sequences. As is the convention in this compendium, a dash (-) indicates concurrence with the top sequence in the alignment; a period (.) indicates a deletion. The carets show where the N-linked glycosylation sites are found in the consensus. The V3 loop is set off from the surrounding sequence by a space on either side to facilitate viewing. Interesting features of the consensus alignment are: 1) Only in the B subtype is GPGR

the most common tip of the V3 loop; globally, GPGQ is more prevalent. 2) A highly conserved N-linked glycosylation site is constitutively absent in the C subtype, proximal to the first cysteine (C) in loop. 3) The D subtype consensus has 34 amino acids from cysteine (C) to cysteine (C) rather than the more common 35; at the point where the deletion occurs, it is not uncommon to find insertions of 2 to 4 amino acids, as can be observed in the sequence alignments. 4) The D subtype has two arginine (R) residues in the V3 loop that are uncharacteristic relative to the other consensus sequences; positively charged amino acids in these positions may result in a syncytia inducing, non-monocytropic phenotype (Fouchier RAM, et al., *J. Virol.* **66**:3183–3187 (1992)). 5) A higher degree of variation is seen in the region just downstream of the V3 loop than within it. This difference is also observed internally among the sequences of the different subtypes. 6) The A, C, G and H consensus sequences have very similar V3 loop sequences.

## V3 Loop Amino Acids

The following pages present amino acid alignments of the V3 loop, arranged by phylogenetic subtype. For each subtype, the number of sequences used to construct the alignment is indicated. The top line in each alignment represents the consensus sequence for that subtype, where consensus simply means the most common amino acid found in each position among the sequences of the given subtype. The subscripts record the frequency with which that amino acid is observed at that location among members of the subtype. An amino acid which is conserved 100% is shown with no subscript. Directly beneath the most common to least common. An asterisk (\*) subscript means less than 0.5% of the sequences had the indicated amino acid at that location. A dash (-) indicates a gap inserted to maintain the alignment. Percentages were rounded to the nearest whole number.

For this year's alignment, the HMMER (version 1.8) hidden Markov model software (Sean Eddy, Dept. of Genetics, Washington U. School of Medicine, St. Louis, MO 63110; eddy@genetics.wustl.edu) was used to objectively align all 1651 sequences. The frequency counts are derived from this alignment. Because each subtypes required different numbers and positions of gaps in order to create the full multiple sequence alignment, some sequences with unusual insertions were trimmed from the HMMER alignment, and a few positions were adjusted by hand, using MASE, prior to printing the full alignment which appears following the country codes description. The sequences which were culled from the alignment after counting frequencies, are appended.

Both the untouched HMMER alignment, and the edited version, will be available via ftp from the LANL HIV database (http://hiv-web.lanl.gov). Questions about these alignments should be directed to (btf@t10.lanl.gov) (505-665-1970).

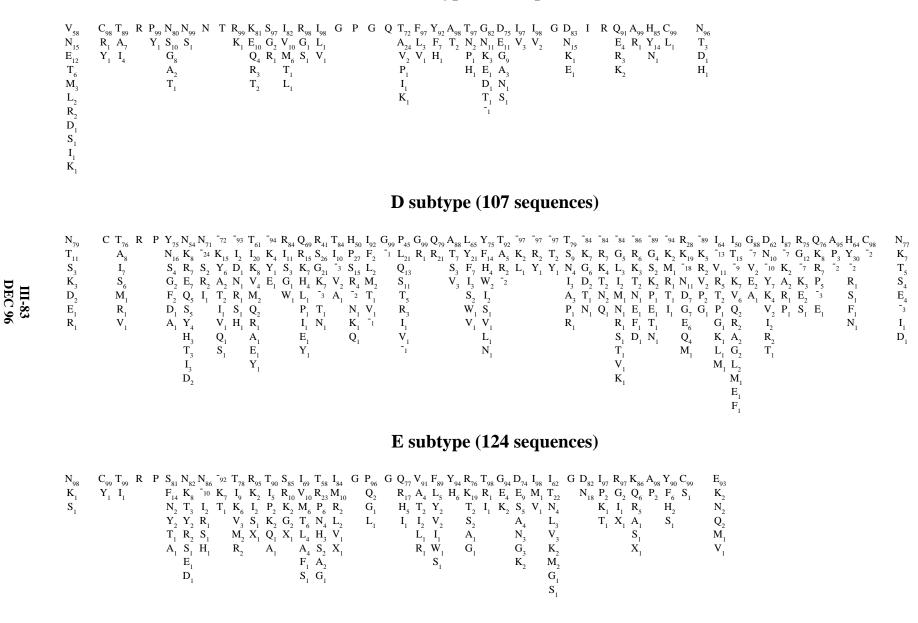
# A subtype (207 sequences)

$N_{_{84}}$	$C T_{76} R P_{99} N_{72} N_{97} N_{$	$V_{97}T_{97}R_{97}K_{67}S_{83}V_{54}R_{51}I_{95}G_{99}P_{98}$	$G Q_{91} A_{65} F_{95} Y_{95} A_{93} T_{82} G_{86} D_{64} I_{95} I_{93} G_{86} D_{64} I_{95} I_{95} G_{86} G_{86$	$G_{99}D_{80}I_{99}R_{98}Q_{81}A_{99}H_{77}C$	N <sub>82</sub>
T <sub>8</sub>	$I_{16} T_* L_1 G_{14} K_1$	$\mathbf{F}_{3} \mathbf{R}_{1} \mathbf{S}_{1} \mathbf{R}_{15} \mathbf{G}_{13} \mathbf{I}_{42} \mathbf{H}_{42} \mathbf{L}_{2} \mathbf{A}_{1} \mathbf{S}_{2}$		$E_* N_{16} T_* K_1 K_{14} V_* Y_{20} $	<b>T</b> <sub>9</sub>
$D_{3}$	$S_{6} = H_{*} S_{10} D_{1}$	$I_{1} N_{*} T_{14} R_{2} M_{1} P_{3} M_{2} R_{*}$	$\mathbf{K}_{2} \mathbf{S}_{3} \mathbf{L}_{1} \mathbf{H}_{1} \mathbf{I}_{2} \mathbf{S}_{4} \mathbf{D}_{5} \mathbf{A}_{8} \mathbf{M}_{1} \mathbf{V}_{1}$	$\mathbf{F}_{1} \mathbf{K}_{1} \mathbf{F}_{1} \mathbf{K}_{2} \mathbf{F}_{1}$	$E_{5}$
$\mathbf{S}_{2}$	$\mathbf{V}_{1}$ $\mathbf{Y}_{1}$ $\mathbf{R}_{*}$	$\mathbf{K}_{*} \mathbf{K}_{*} \mathbf{Q}_{2} \mathbf{N}_{1} \mathbf{L}_{*} \mathbf{S}_{2} \mathbf{F}_{1}$	$\mathbf{G}_{*} \mathbf{V}_{2} \mathbf{S}_{1} \mathbf{I}_{*} \mathbf{V}_{*} \mathbf{A}_{3} \mathbf{N}_{*} \mathbf{G}_{7} \mathbf{I}_{1} \mathbf{I}_{1}$	$S_1 \xrightarrow{-} E_* N_*$	$D_1$
<b>K</b> <sub>1</sub>	$M_*$ $A_*$ $T_*$	$V_{*} L_{*} E_{1} M_{*} D_{*} N_{1} T_{*}$	$\mathbf{S}_* \ \mathbf{P}_* \ \mathbf{W}_* \ \mathbf{C}_* \ \mathbf{G}_* \ \mathbf{C}_2 \qquad \mathbf{N}_4 \qquad \mathbf{G}_*$	$P_*$ $L_*$ $S_*$	$I_1$
$\mathbf{I}_{1}$	$A_*$ $T_*$	$M_* E_* A_* C_*$	$I_*$ $I_1$ $R_2$	$E_*$ - $C_*$	<b>K</b> <sub>1</sub>
$\mathbf{F}_*$	$\mathbf{X}_{*}$	$\mathbf{S}_* = \mathbf{W}_* \mathbf{K}_*$	$\mathbf{N}_{1}$ $\mathbf{K}_{1}$	* *	$V_*$
$\mathbf{E}_*$	$\mathbf{F}_{*}$	$A_*$ $ G_*$	K <sub>*</sub> <sup>-1</sup>		-*
		$\mathbf{X}_{*}$	$\mathbf{P}_*$ $\mathbf{S}_1$		
			$Q_*$		

# **B** subtype (975 sequences)

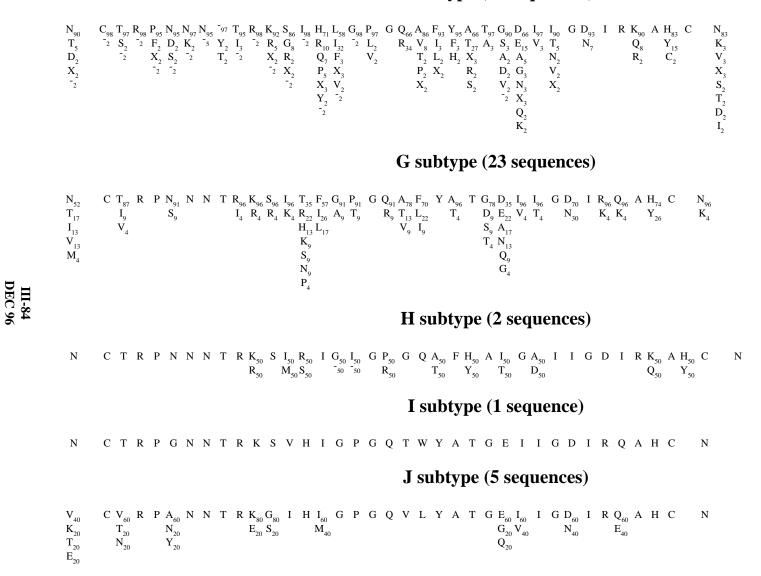
N <sub>91</sub>	$C T_{90} R_{99} P_{99}$	N <sub>86</sub> N <sub>98</sub> N	95 96 7	$T_{97} R_{97}$	<sub>95</sub> K <sub>84</sub>	S <sub>74</sub>	I <sub>96</sub> H	57 I 68	G <sub>97</sub>	P <sub>92</sub>	G <sub>99</sub>	R <sub>77</sub>	A <sub>88</sub>	F <sub>71</sub>	Y <sub>88</sub>	T <sub>57</sub>	T <sub>95</sub>	G <sub>88</sub>	E <sub>36</sub>	I <sub>96</sub>	I <sub>91</sub> (	G <sub>99</sub> I	) 87	$I_{98}$	R <sub>99</sub>	Q <sub>87</sub>	Α	$H_{91}$	С	$N_{90}$
T <sub>3</sub>	$S_* I_8 S_* L_1$	$S_7 T_1$	4 Y <sub>1</sub>	I <sub>1</sub> S	$_{2} R_{15}$	G <sub>18</sub>	V <sub>3</sub> P	$_{17} L_{19}$	A2	$W_4$	$\mathbf{R}_*$	$Q_9$	V <sub>5</sub>	W <sub>17</sub>	$F_4$	A440	$R_1$	7	D <sub>24</sub>	V <sub>2</sub>	$V_4$	E <sub>1</sub> 1	N <sub>11</sub>	$T_1$	<b>K</b> <sub>1</sub>	$K_8$	$\mathbf{P}_{*}$	Y <sub>8</sub> '	*	-4
H <sub>2</sub>	$R_* V_1 T_*$	$G_3 K_1 H$	* K <sub>1</sub>	K <sub>1</sub> I	X.	R <sub>7</sub>	L <sub>1</sub> N	$I_8 M_1$	$_{1}E_{*}$	$L_2$	$\mathbf{K}_*$	$K_8$	T <sub>5</sub>	$L_5$	$H_4$	2	$A_1$	$E_2$	Q <sub>19</sub>	-1	T <sub>3</sub>	*	E*	$V_*$	$\mathbf{S}_*$	R <sub>3</sub>	$\mathbf{S}_*$	$Q_1$ S	S <sub>*</sub>	T <sub>3</sub>
D <sub>1</sub>	$X_* S_1 X_*$	H <sub>1</sub> D <sub>*</sub> S	* T <sub>1</sub>	V <sub>*</sub> K	T.	$H_*$	М <sub>*</sub> Т	$V_{1}$	$Q_{*}$	$Q_1$	$\mathbf{E}_*$	<b>S</b> <sub>3</sub>	S <sub>1</sub>	V.,	$V_1$	$G_*$	1	<b>R</b> <sub>1</sub>	R <sub>5</sub>	$T_*$	A <sub>1</sub>	V <sub>*</sub> (	$G_*$	$M_{*}$	*	E <sub>1</sub>	$T_*$	$\mathbf{F}_{*}$		I
K <sub>1</sub>	E <sub>*</sub> -*	T <sub>1</sub> I <sub>*</sub> I	* E*	$S_* X$	* Q*	$\mathbf{D}_*$	K <sub>*</sub> S	5 F.	$R_*$	$\mathbf{F}_{1}$	$W_*$	$G_2$	$\mathbf{R}_*$	$I_2$	$R_1$	$V_*$	$I_1$	K <sub>1</sub>	$G_4$	$L_*$	$L_*$	R <sub>*</sub> (	$Q_*$	$\mathbf{K}_*$	$G_{*}$	$L_1$		$A_*$		$\mathbf{D}_*$
<b>S</b> <sub>1</sub>	$\mathbf{A}_{*}$	$Y_1 R_* X$	• F.	$A_* M$	$I_* E_*$	$C_*$	T <sub>*</sub> Y	, T.		$G_1$		-*	*	$\mathbf{Y}_1$	$I_1$	$\mathbf{R}_*$	$N_{*}$	$\mathrm{D}_*$	$K_4$	$K_*$	$X_*$		*	$L_*$	$T_*$	$\mathrm{H}_{*}$		$\mathbf{R}_*$		$\mathbf{S}_*$
$E_*$	$\mathbf{X}_{*}$	D <sub>*</sub> H <sub>*</sub>	$Q_*$	R <sub>*</sub> G	* N*	$T_*$	F <sub>*</sub> R	έ <sub>2</sub> Κ <sub>*</sub>		$\mathbf{S}_*$		$\mathbf{E}_*$	$\mathbf{X}_*$	$S_*$	N <sub>*</sub>	$\mathbf{S}_*$	$\mathbf{S}_*$	$A_*$	A,	$\mathbf{R}_*$	-*	]	K*	-*	$\mathbf{I}_*$	$X_*$		$X_*$		$\mathbf{E}_{*}$
$I_*$	$L_*$	$X_* Y_*$	$\mathbf{S}_{*}$	X <sub>*</sub> Q	* I*	$\mathbf{X}_*$	Y <sub>*</sub> Q	$\bar{\mathbf{p}}_1 \mathbf{S}_*$		$A_*$		$N_*$	$G_{*}$	$\mathbf{M}_{*}$	$W_*$	$\mathbf{X}_*$	$\mathbf{K}_*$	$T_*$	N <sub>2</sub>	$\mathbf{E}_{*}$	$K_*$		$A_*$			$A_{*}$		$N_*$		$\mathbf{K}_{*}$
-*	$M_{*}$	$F_*$		-* A	*	$\mathbf{V}_{*}$	⁻∗ A	Υ <sub>1</sub> Υ <sub>*</sub>		$V_*$		$X_*$	$\boldsymbol{Q}_{\ast}$	$C_*$	$\mathrm{D}_*$	$\mathbf{K}_*$	$\mathbf{X}_*$	$X_*$	$H_1$		$F_*$	]	$H_*$			-*				$V_*$
Y <sub>*</sub>		K <sub>*</sub>		L	*	$N_*$	G	; ;		$M_*$		$A_*$		$\mathbf{X}_*$	$T_*$		$\mathbf{Y}_{*}$	$\mathbf{S}_{*}$	1		$M_*$		Τ.							$X_*$
$X_*$				Т	*		Х	*		$R_*$				$\mathbf{P}_{*}$	$C_*$		$\mathbf{P}_{*}$		$\mathbf{S}_*$		$Q_*$		V.							$\mathbf{Y}_{*}$
				Ε	*		Ι	*		$T_*$					$L_*$				$X_*$				Y <sub>*</sub>							$H_{*}$
							F	*							$\mathbf{S}_*$				$V_*$				I.							
							K	*											$\mathbf{P}_*$											
							-	*											$\mathbf{Y}_{*}$											
																			$T_*$											

## C subtype (119 sequences)

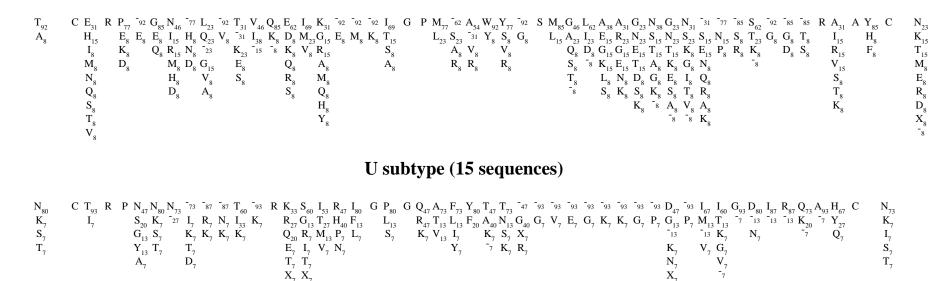


HIV-1 V3 Region

### **F** subtype (59 sequences)



## O subtype (13 sequences)



É,

III-85 DEC 96

HIV-1 V3 Region

### **V3** Loop Variation

**Summary of variations in the tetrameric tip of the V3 loop.** This table is a tally of the different tetramers observed in the 1651 individuals analyzed. This tip is thought to form a turn, and is the focal point of the potent neutralizing antibody epitopes that have been mapped to the V3 loop, as well as of T cell epitopes. Each column shows the number of occurrences of a given tetramer in either the entire 976 sequences (combined), or in subsets consisting of subtypes A–O, and the unclassified sequences (U). Underneath the column heading is the number of sequences in each category. The most common form found in each subtype is highlighted in bold lettering. In the B subtype, GPGR is the predominant form, however globally GPGQ is more common.

	Combined	А	В	С	D	Е	F	G	Η	Ι	J	0	U
Totals	1651	208	975	119	107	124	59	23	2	1	5	13	15
GPGR	711	10	643		17	18	17						6
GPGQ	590	184	91	119	31	95	39	19	1	1	5		5
GPGK	81	4	76										1
GWGR	34		34										
GPGS	26	1	25										
GPGG	22	2	20										
GLGQ	21				20								1
GLGR	19		15		1	1	1						1
APGR	17	1	16										
GSGQ	16	3			12								1
GQGQ	13		1		12								
GQGR	9		6		1	2							
GPMA	8											8	
GPGH	6					6							
GFGR	6		6										
GTGQ	5				5								
GRGQ	5	1			3				1				
GVGR	4		2		1		1						
GSGR	4	1	3										
GPRR	3		3										
GPGA	3		3										
EPGR	3		3										
APGS	3		3										
APGQ	3	1						2					
GTGR	2							2					
GPKR	2		2										
GMGR	2		2										
GGGR	2		2										
GAGR	2		2										
AGGR	2		2										
GLGS	1		1										
GTGG	1		1										
GPMR	1											1	
GLRQ	1				1								
AQGR	1				1							4	
GPMS	1											1	
GPLR	1		1									1	
GARR	1		1			1							
GGGQ	1		4			1							
GPWG	1		1										
GIGQ	1		1		1								
RPGR	1		1										

	Combined	А	В	С	D	Е	F	G	Н	Ι	J	0	U
Totals	1651	208	975	119	107	124	59	23	2	1	5	13	15
GRGR	1		1										
GPGR	1		1										
GQGI	1					1							
GPLS	1											1	
GPWG	1		1										
GPGN	1		1										
GPGX	1		1										
GPGE	1		1										
GPEK	1		1										
GLGK	1		1										
GARR	1		1										
AWGR	1		1										
APGG	1		1										
AGGK	1		1										
AQGR	1				1								
GIGQ	1				1								
GGRA	1				1								
*PGR	1						1						

**Summary of variations in the octameric tip of the V3 loop.** This table is a tally of the different octamers observed in the 1651 individuals analyzed. This table is structured the same as the tetramer table on the previous pages. Amino acid changes proximal to the tip can influence the specificity of anti-V3 neutralizing antibodies. The forms that were found only once in the data set are not shown here, to save space, and are summarized in a row labeled "unique."

	Combined	Α	В	С	D	Е	F	G	Н	Ι	J	0	U
Totals	1651	208	975	119	107	124	59	23	2	1	5	13	15
HIGPGRAF	279	3	269		2		3						2
RIGPGQTF	136	46		82	1	2	1	3					1
RIGPGQAF	75	45		27		1		1					1
PIGPGRAF	71	1	69				1						
HIGPGQAF	62	56	2				4						
NIGPGRAF	59		58										1
HLGPGQAW	44		44										
TIGPGQVF	39					39							
PLGPGQAW	31		31										
HLGPGQAF	31	2	1				28						
SIGPGRAF	26	1	25										
HIGPGKAF	25	2	23										
RIGPGQVF	23	3		2		16	1						1
TIGPGRAF	19		19										
HIGPGQAL	16		2		13			1					
YIGPGRAF	15		15										
PIGLGQAL	14				13								1
HMGWGRAF	14		14										
HIAPGRAF	14	1	13										
PLGPGRAW	11		11										
HLGPGRAW	10		10										
HIGPGSAF	10		10										
HIGPGRAY	10				10								
TMGPGQVF	9					9							
PMGPGRAF	9		9										
HMGPGRAF	9	1	8										
RIGPGRVF	8					8							
PIGPGQAF	7	3			1	2		1					
HMGWGRTF	7		7										
HIGPGRTF	7		7										
TIGPGRVF	6		1			5							
QIGPGRAF	6		5				1						
PMGPGKAF	6		6										
PLGPGKAW	6		6										
NIGPGRAW	6		6										
HLGWGRAF	6		6										
HIGPGRAV	6		6										
YLGPGRAF	5		5										
RFGPGQAF	5	1	-				1	1					2
PIGPGKAF	5	-	5				-	-					_
NMGPGRAF	5		5										

	Combined	А	В	С	D	Е	F	G	Η	Ι	J	0	U
Totals	1651	208	975	119	107	124	59	23	2	1	5	13	15
NIGPGQVF	5					5							
HLGPGRAF	5	1	1				2						1
HLGPGGAF	5		5										
HIGSGQAL	5				5								
HIGPGRAW	5		5										
HIGPGRAL	5		2		1		1						1
HIGPGRAI	5		5										
HIGPGGAF	5		5										
SIGQGQAL	4				4								
SIGPGQAF	4	3							1				
~ PIGPGRAW	4		4										
PIGPGQVF	4					4							
KIGPMAWY	4											4	
HMGPGKAF	4		4										
HMGLGRAF	4		4										
HLGPGQAL	4		2		2								
HIGPGRVF	4		3		-		1						
HIGPGRSF	4	2	2				1						
HIGPGQVF	4	1	-			3							
HIGPGQTF	4	4				5							
HIGPGQAI	4	2			2								
YIGPGRAV	3	2	3		2								
YIGPGRAS	3		3										
TMGPGRAS	3		3										
	3		3										
TMGPGRVW	3		3										
TLGPGRVY	3		1			2							
TIGPGRVY	3		1			2		3					
TFGPGQAF	3		3					5					
SLGPGRAW	3		3										
SIGPGRAW	3		5	2	1								
RIGPGQTL	3	3		2	1								
RIGPGQSF	3	5	1				2						
PLGPGRAF	3	2	1		1		2						
NIGPGQAF	3	1	2		1								
HMGPGKTF	3	1	2 3										
HLGQGRAW	3		5		2								
HIGTGQAL	3				3 3								
HIGSGQAY	3 3				3						3		
HIGPGQVL			2								3		
HIGPGQAW	3	n	3	1									
GIGPGQTF	3	2		1		2							
AIGPGQVF	3		2			3							
XIGPGRAF	2		2			2							
TRGPGHVF	2		2			2							
TMGPGRVL	2		2										
TMGPGKVF	2		2			~							
TMGPGHVF	2					2		-					
TLGPGQAF	2					_		2					
TIGPGQVL	2					2							

## V3 Loop Variation

	Combined	А	В	С	D	Е	F	G	Η	Ι	J	0	U
Totals	1651	208	975	119	107	124	59	23	2	1	5	13	15
TIGPGQIF	2					2							
SMGPGRAF	2		2										
SLGPGKAW	2		2										
SIGQGRVL	2					2							
SIGQGQTL	2				2								
SIGPGRVW	2		2										
SIGLGQAL	2				2								
SFGPGQAF	2							2					
RIGPMAWY	2											2	
RIGPGSAF	2		2										
RIGPGRVI	2						2						
RIGPGRTF	2		1										1
RIGPGRAV	2		2										
RIGPGRAF	2		2										
RIGPGQAL	2				2								
~ QLGPGRAW	2		2										
~ PLGPGRVW	2		2										
PIGSGQAL	2				2								
PIGRGQAL	2				2								
~ PIGLGQAY	2				2								
~ PIAPGSAW	2		2										
KIGPGQTF	2	1			1								
~ KFGTGRVL	2							2					
HVGPGQAF	2				1		1						
~ HMGPGRAL	2				2								
HMGPGQVL	2										2		
HMGPGGAF	2		2										
HLGPGKAW	2		2										
HLGPGKAF	2		2										
HLGLGRAF	2		2										
HLGFGRAL	2		2										
HIGSGRAF	2	1	1										
HIGSGQAI	2	-	_		1								1
HIGPGSAL	2		2		-								-
HIGPGQAY	2		-		2								
HIGPGKVF	2		2		-								
HIGLGRAY	2		-		1								1
HIGGGRTL	2		2										•
HIEPGRAF	2		2										
HIEPGRAF	2		4		1			1					
GIGPGQAL	2		2		1			1					
AIGPGRIV	2		2										
UNIQUE	189	26	93	3	26	8	10	11	2	3	7		